

Score Function

虎尾科技大學
40723115
林于哲

January 16 2021

Contents

1	Log Derivative Trick	1
2	Score Functions	1
3	Score Function Estimators	2

1 Log Derivative Trick

機器學習涉及操縱機率。這個機率通常包含 normalised-probabilities 或 log-probabilities。能加強解決現代機器學習問題的關鍵點，是能夠巧妙的在這兩種型式間交替使用，而對數導數技巧就能夠幫助我們做到這點，也就是運用對數導數的性質。

2 Score Functions

對數導數技巧的應用規則是基於參數 θ 梯度的對數函數 $p(x : \theta)$ ，如下：

$$\nabla_{\theta} \log p(x : \theta) = \frac{\nabla_{\theta} p(x : \theta)}{p(x : \theta)}$$

$p(x : \theta)$ 是 likelihood ; function 參數 θ 的函數，它提供隨機變量 x 的概率。在此特例中， $\nabla_{\theta} \log p(x : \theta)$ 被稱為 Score Function，而上述方程式右邊為 score ratio(得分比)。

The score function has a number of useful properties:

- The central computation for maximum likelihood estimation. Maximum likelihood is one of the dominant learning principles used in machine learning, used in generalised linear regression, deep learning, kernel machines, dimensionality reduction, and tensor decompositions, amongst many others, and the score appears in all these problems.
- The expected value of the score is zero. Our first use of the log-derivative trick will be to show this.

$$\begin{aligned}\mathbb{E}_{p(\mathbf{x};\theta)}[\nabla_{\theta} \log p(\mathbf{x};\theta)] &= \mathbb{E}_{p(\mathbf{x};\theta)} \left[\frac{\nabla_{\theta} p(\mathbf{x};\theta)}{p(\mathbf{x};\theta)} \right] \\ &= \int p(\mathbf{x};\theta) \frac{\nabla_{\theta} p(\mathbf{x};\theta)}{p(\mathbf{x};\theta)} d\mathbf{x} = \nabla_{\theta} \int p(\mathbf{x};\theta) d\mathbf{x} = \nabla_{\theta} 1 = 0\end{aligned}$$

In the first line we applied the log derivative trick and in the second line we exchanged the order of differentiation and integration. This identity is the type of probabilistic flexibility we seek: it allows us to subtract any term from the score that has zero expectation, and this modification will leave the expected score unaffected (see control variates later).

- The variance of the score is the Fisher information and is used to determine the Cramer-Rao lower bound.

$$\mathbb{V}[\nabla_{\theta} \log p(\mathbf{x};\theta)] = \mathcal{I}(\theta) = \mathbb{E}_{p(\mathbf{x};\theta)}[\nabla_{\theta} \log p(\mathbf{x};\theta) \nabla_{\theta} \log p(\mathbf{x};\theta)^{\top}]$$

We can now leap in a single bound from gradients of a log-probability to gradients of a probability, and back. But the villain of today's post is the troublesome expectation-gradient of Trick 4, re-emerged. We can use our new-found power—the score function—to develop yet another clever estimator for this class of problems.

3 Score Function Estimators